

# Synthetic Health Data: The Good, the Bad, and the Ugly

Bradley Malin, Ph.D.

Accenture Professor of Biomedical Informatics, Biostatistics, & Computer Science

Co-Director, Center for Health Data Science

Co-Director, Center for Genetic Privacy & Identity in Community Settings

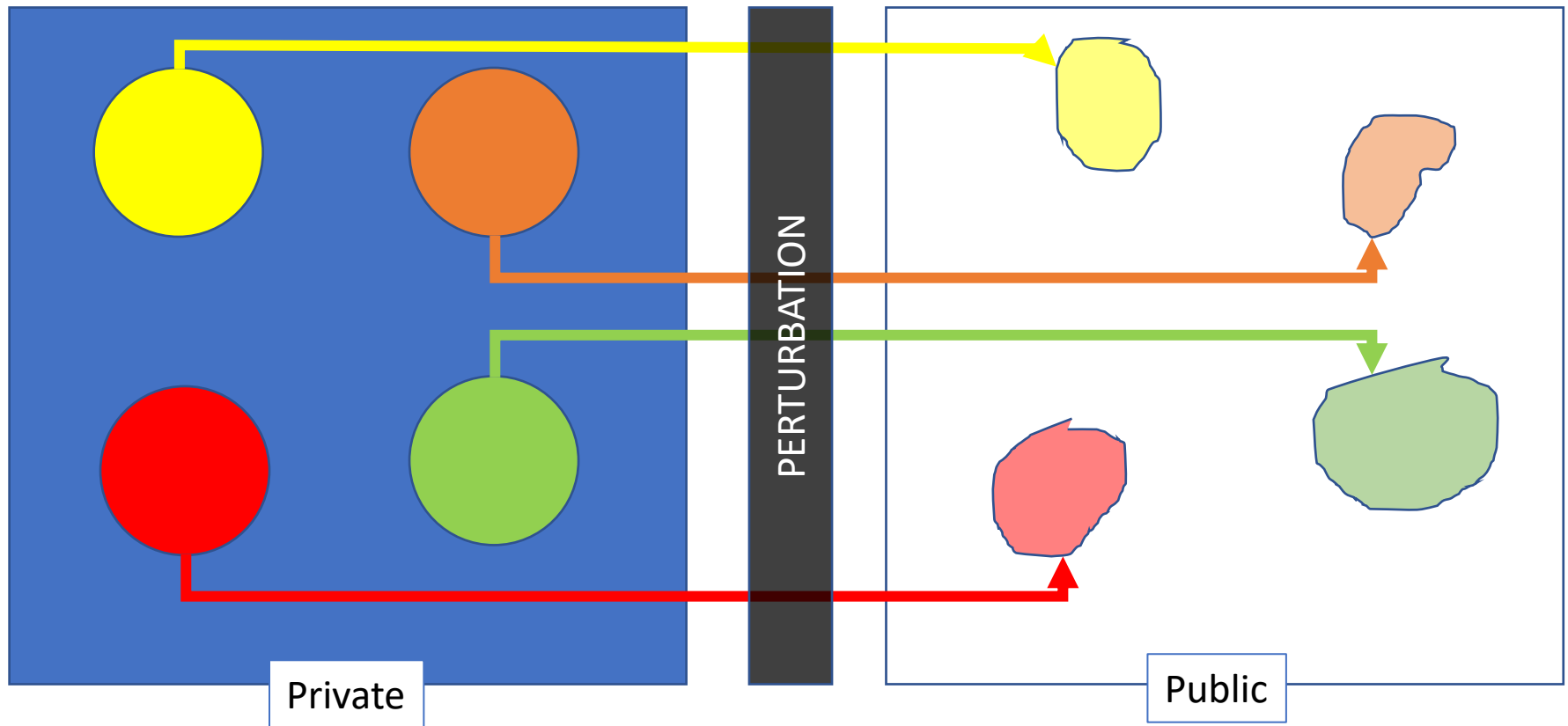
Co-Director, Big Biomedical Data Science Training Program

VANDERBILT  UNIVERSITY

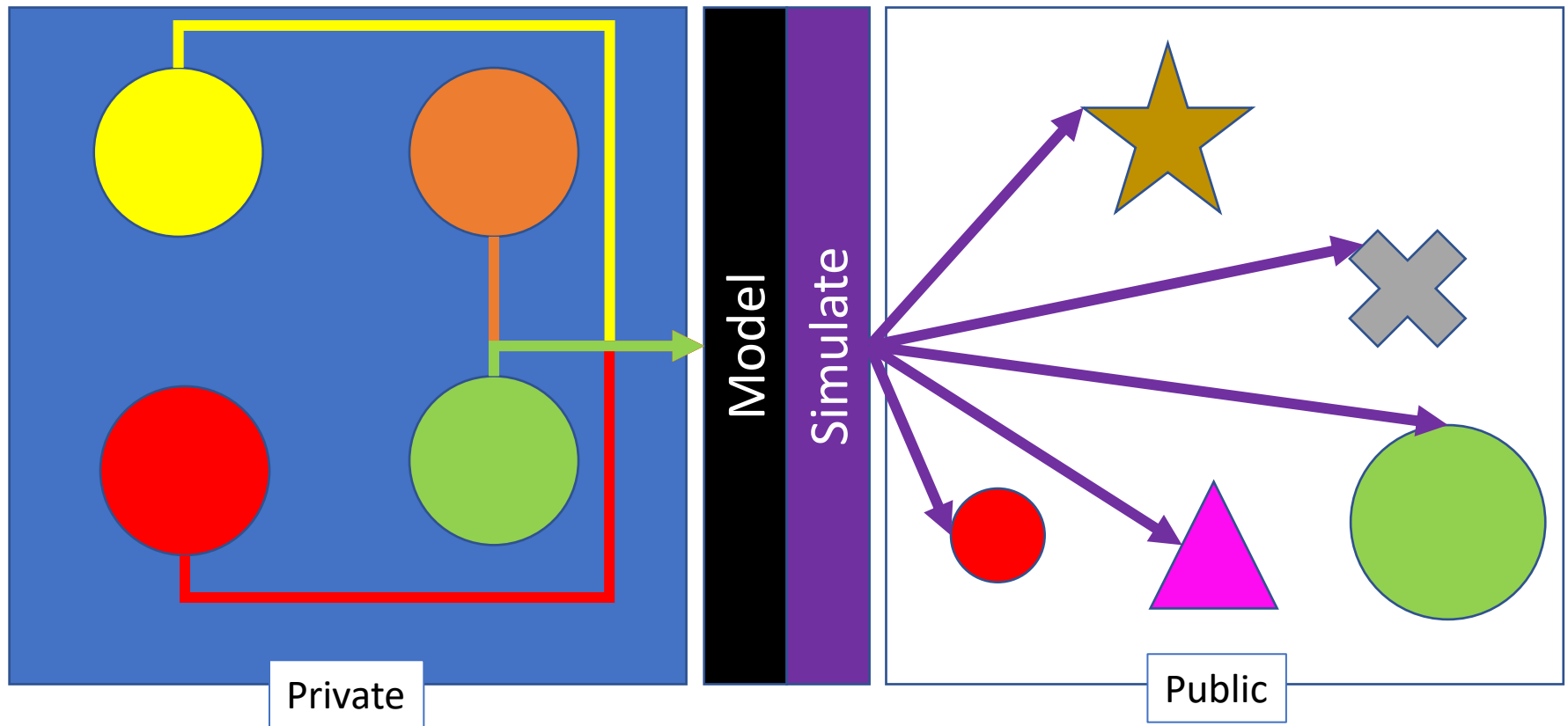
MEDICAL CENTER

November 10, 2021

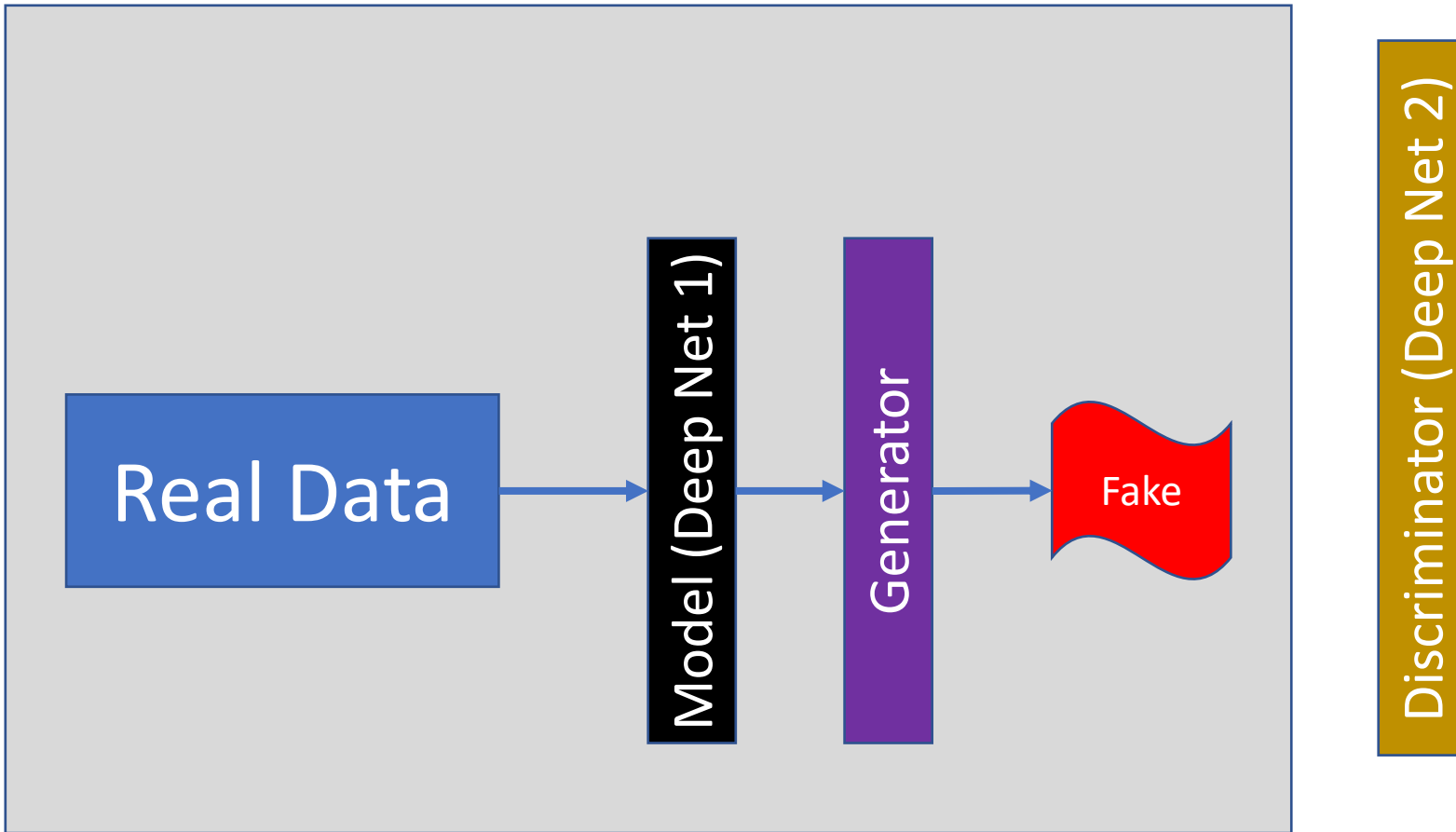
# Ways to Generate Synthetic Data: Perturbation



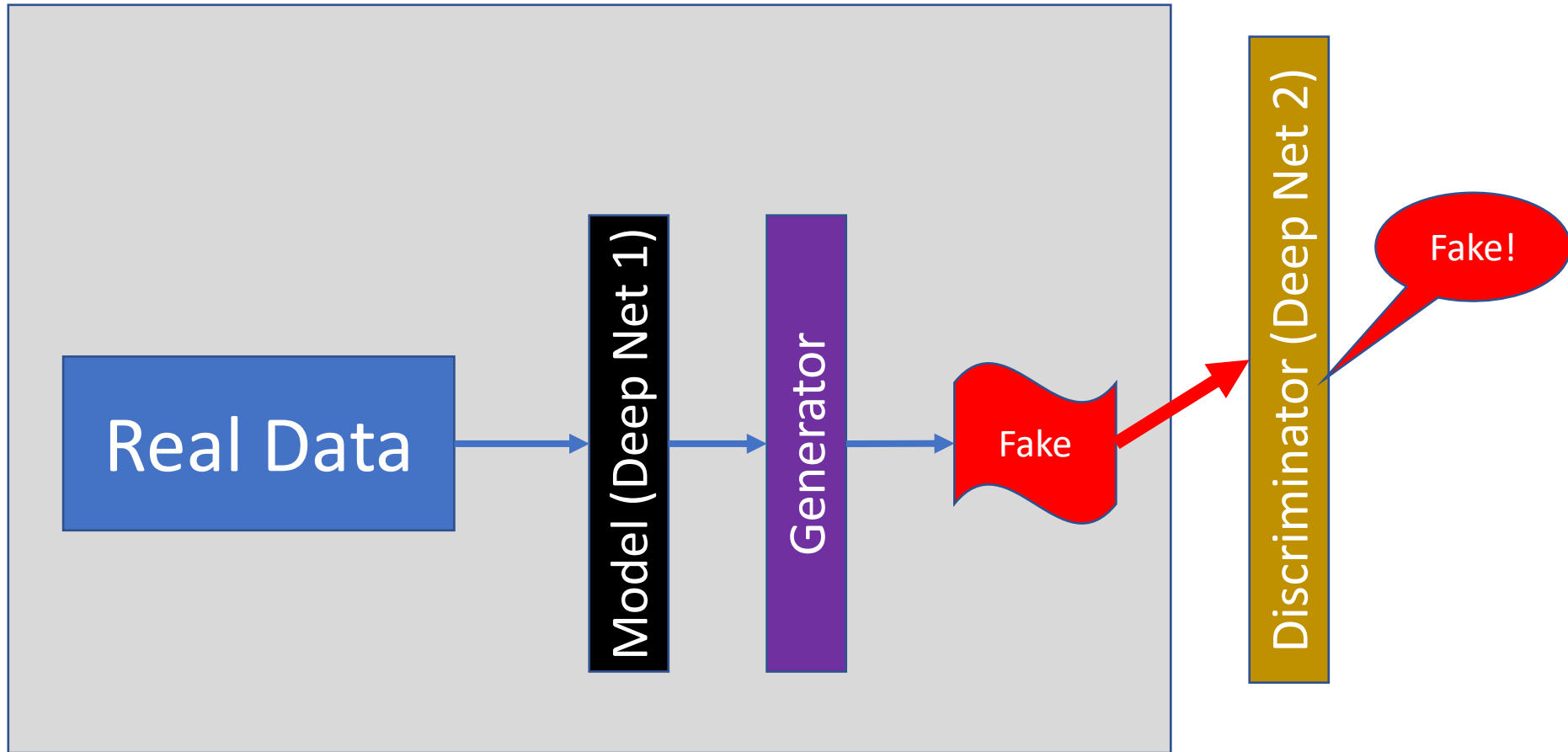
# Ways to Generate Synthetic Data: Simulation



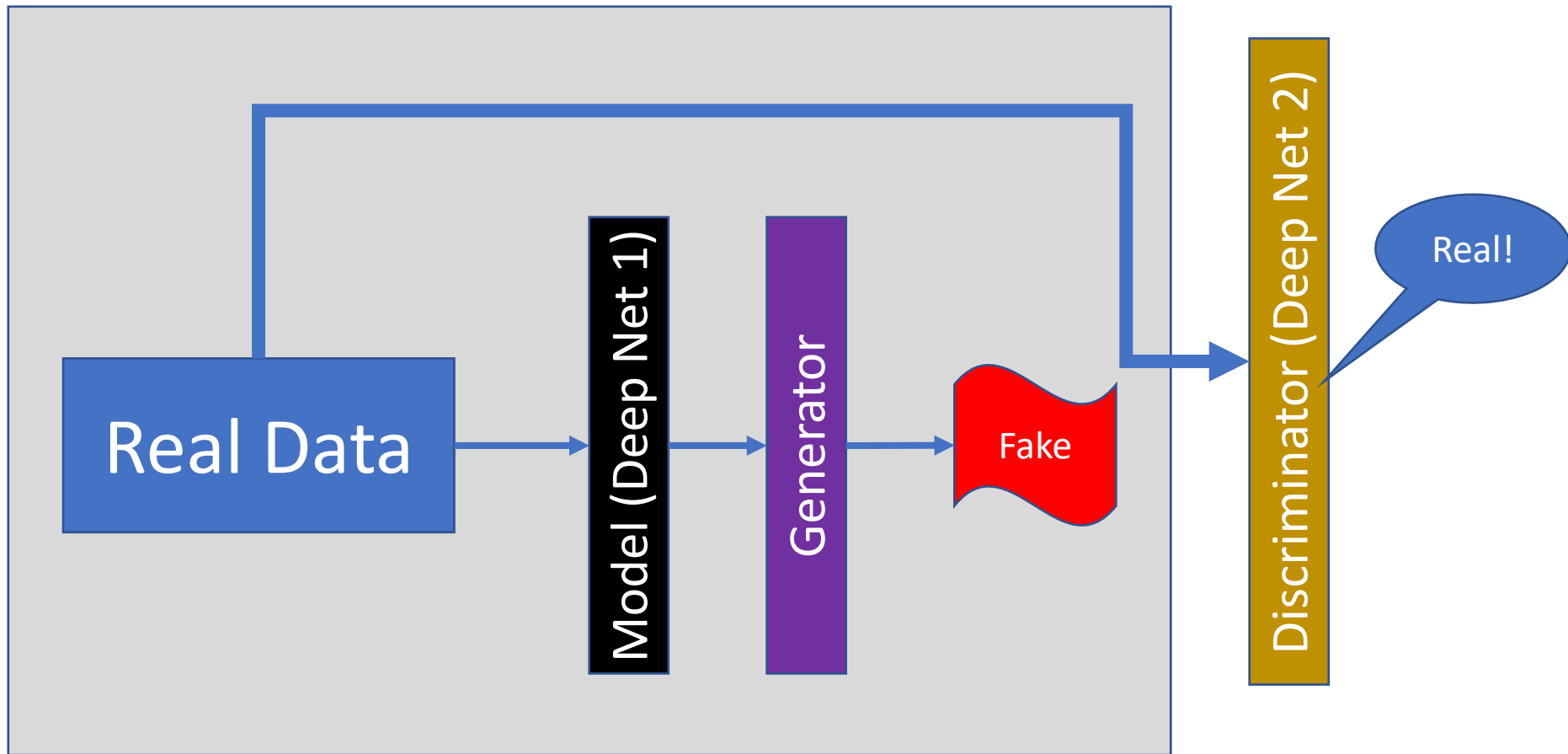
# Generative Adversarial Networks: GANs



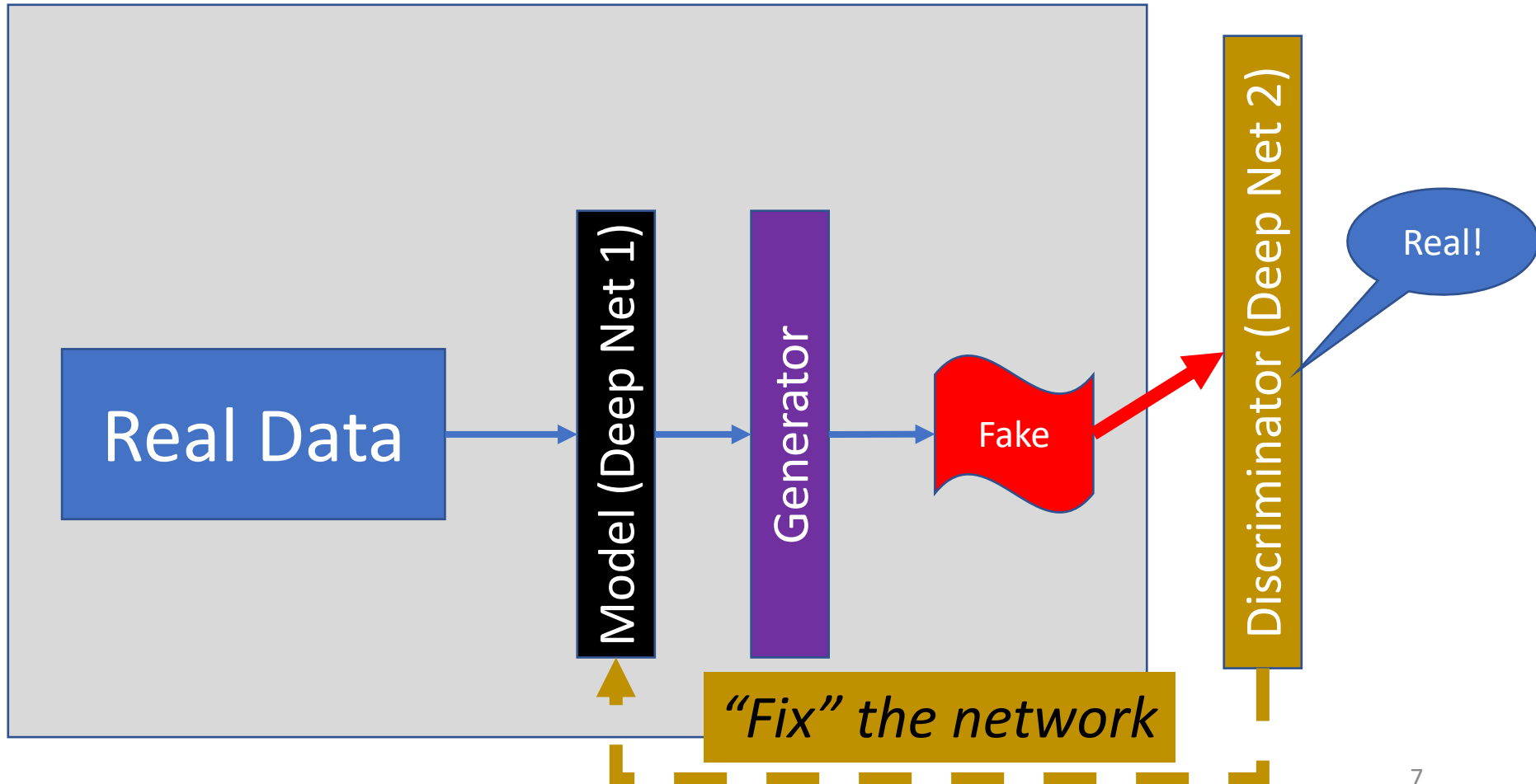
# Generative Adversarial Networks: GANs



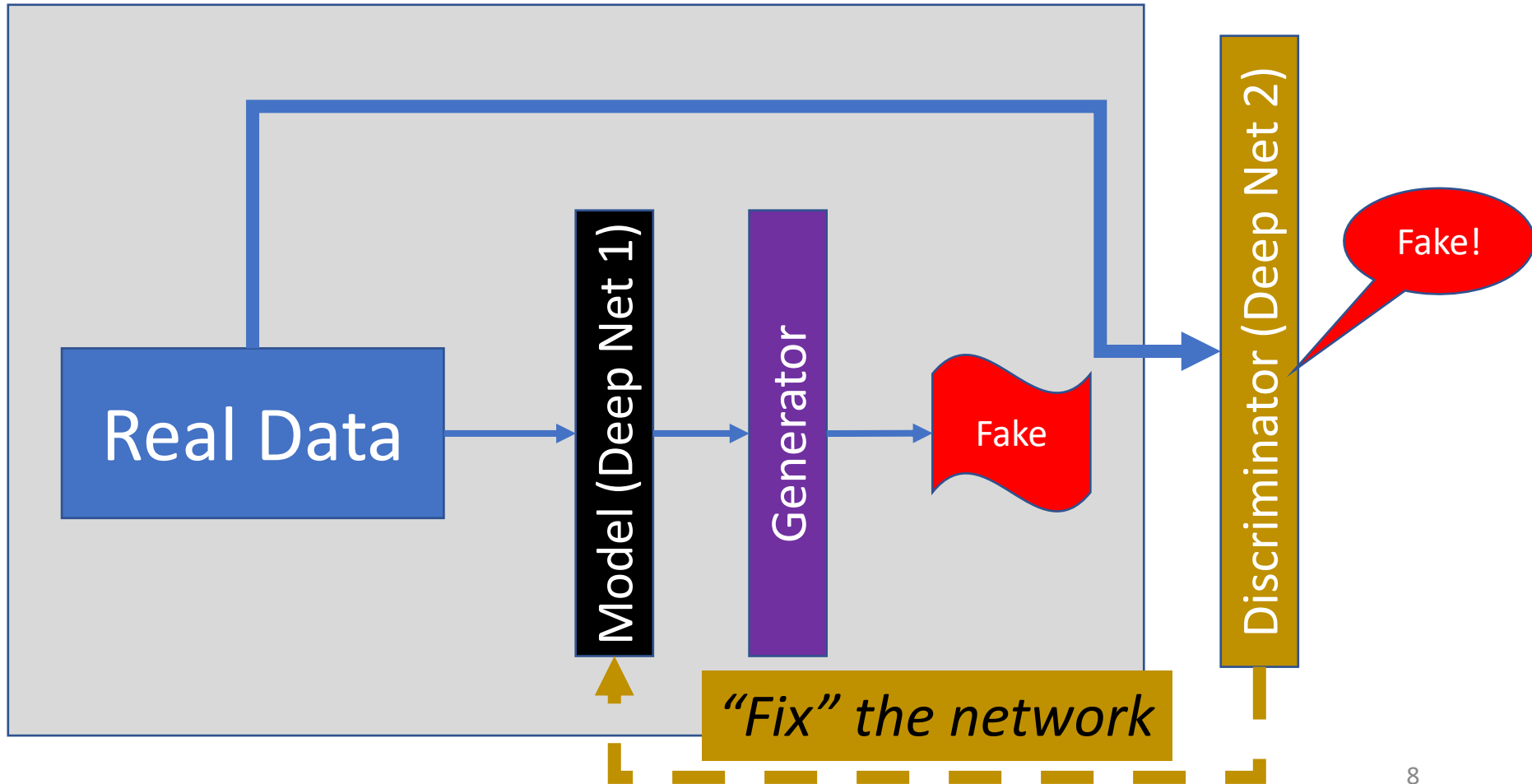
# Generative Adversarial Networks: GANs



# Playing the GAN Game



# Playing the GAN Game





# This is Not a New Principle



**Ian Goodfellow**  
@goodfellow\_ian



4.5 years of GAN progress on face generation.

[arxiv.org/abs/1406.2661](https://arxiv.org/abs/1406.2661) [arxiv.org/abs/1511.06434](https://arxiv.org/abs/1511.06434)

[arxiv.org/abs/1606.07536](https://arxiv.org/abs/1606.07536) [arxiv.org/abs/1710.10196](https://arxiv.org/abs/1710.10196)

[arxiv.org/abs/1812.04948](https://arxiv.org/abs/1812.04948)



## Satisfying Disclosure Restrictions With Synthetic Data Sets

Jerome P. Reiter<sup>1</sup>

To avoid disclosures, Rubin proposed so that (i) no unit in the released data and (ii) statistical procedures that are In this article, I show through simulation from synthetic data in a variety of scenarios proportional to size sampling, two-stage provide guidance on specifying the methods the benefit of including design variables

*Key words:* Confidentiality; disclosure

The screenshot shows the website for JPrivacy Confidentiality. The navigation bar includes links for Current, Archives, Announcements, TPDP workshop, and Submissions. The breadcrumb trail is Home / Archives / Vol. 1 No. 1 (2009): Inaugural Issue / Articles. The article title is 'Estimating Risks of Identification Disclosure in Partially Synthetic Data'. The authors listed are Jerome P. Reiter (Department of Statistical Science, Duke University, Durham, NC) and Robin Mitra (University of Southampton, Southampton). The article is published as a PDF on April 1, 2009, with a DOI of https://doi.org/10.29012/jpc.v1i1.567. The keywords are Confidentiality, Public use data, and To limit disclosures, statistical agencies.

The screenshot shows the article 'Inferentially Valid, Partially Synthetic Data: Generating from Posterior Predictive Distributions not Necessary' from the Journal of Official Statistics, Vol. 28, No. 4, 2012, pp. 583–590. The authors are Jerome P. Reiter<sup>1</sup> and Satkartar K. Kinney<sup>2</sup>. The abstract states: 'To avoid disclosures in public use microdata, one approach is to release partially synthetic data sets. These comprise the units originally surveyed with some collected values, for example sensitive values at high risk of disclosure or values of key identifiers, replaced with multiple imputations. In practice, partially synthetic data typically are generated from Bayesian posterior predictive distributions; that is, one draws repeated values of parameters in the synthesis models before generating data from them. We show, however, that inferentially valid, partially synthetic data can be generated by fixing the parameters of the synthesis models at their modes. We do so with both a theoretical example and illustrative simulation studies. We also discuss implications of these results for agencies generating synthetic data.' The key words are Confidentiality; disclosure; imputation; microdata; privacy; survey.

# This is Not a New Principle

(Choi MLHC 2017)

Proceedings of Machine Learning for Healthcare 2017

JMLR W&C Track Volume 68

## Generating Multi-label Discrete Patient Records using Generative Adversarial Networks

Edward Choi<sup>1</sup>

Siddharth Biswal<sup>1</sup>

Bradley Malin<sup>2</sup>

Jon Duke<sup>1</sup>

Walter F. Stewart<sup>3</sup>

Jimeng Sun<sup>1</sup>

MP2893@GATECH.EDU

SBISWAL7@GATECH.EDU

BRADLEY.MALIN@VANDERBILT.EDU

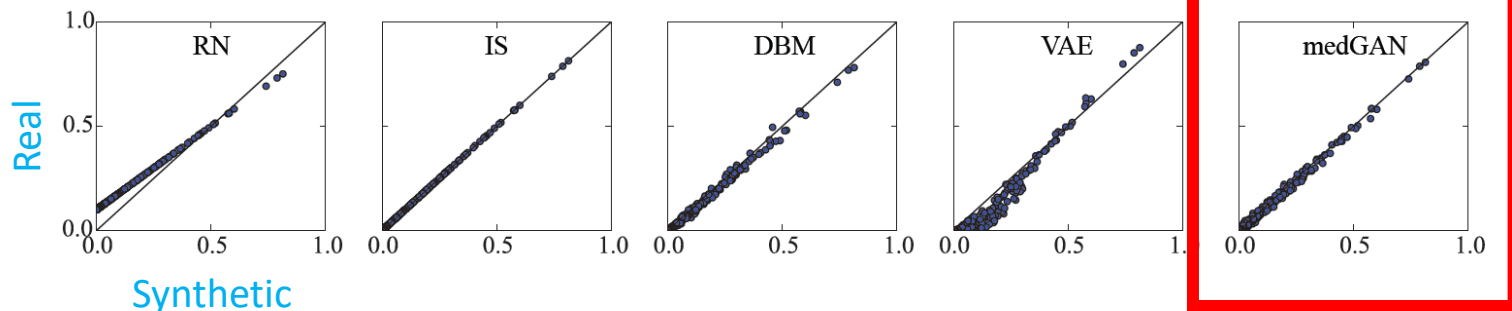
JON.DUKE@GATECH.EDU

STEWARWF@SUTTERHEALTH.ORG

JSUN@CC.GATECH.EDU

<sup>1</sup>GEORGIA INSTITUTE OF TECHNOLOGY   <sup>2</sup>VANDERBILT UNIVERSITY   <sup>3</sup>SUTTER HEALTH

- Sutter Health & MIMIC
- Demographics, Diagnoses, Procedures, & Meds
- Prediction of presence / absence clinical concept



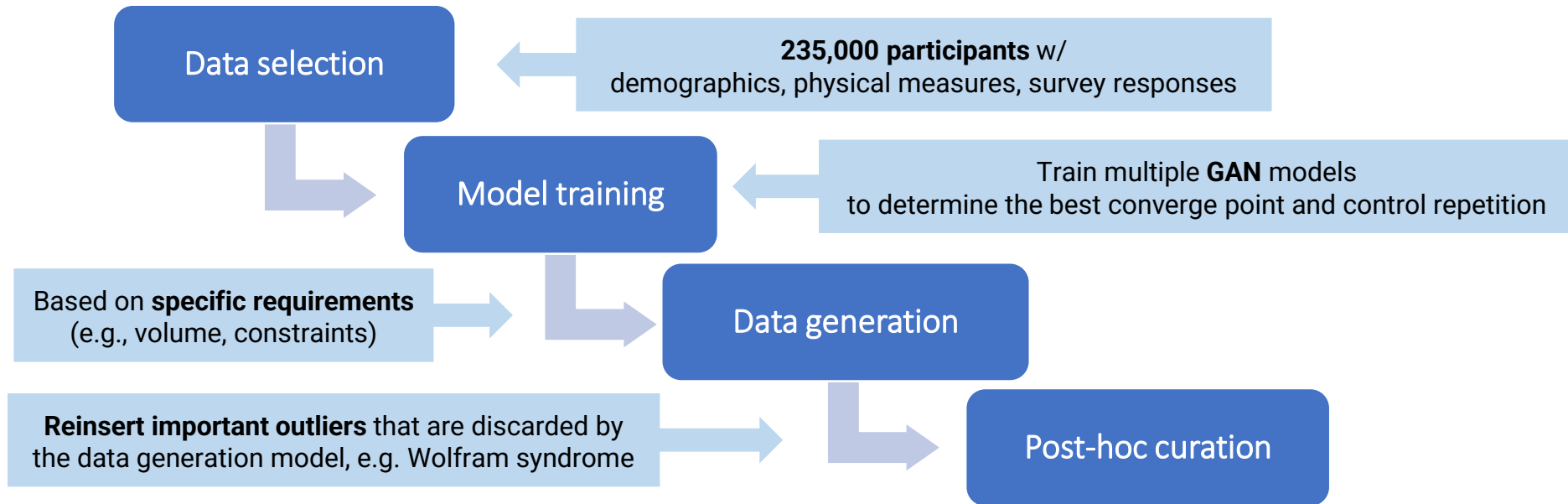
# Limitations

- Autoencoder induced noise and hurt learning
- Evaluation measures based on superficial aspects of data gave false impression of merits of simulation
- Focus on all EHR data led to overrepresentation of common associations

# Evolution

- Better training (Wasserstein distance) and evaluation methods (latent dimensions) (Zhang JAMIA 2020)
- Enabling constraints (e.g., preventing women from having prostate cancer) (Yan AMIA 2020)
- Move from static to longitudinal data: think LSTMs + GANs (Zhang JAMIA 2021)

# Building a Synthetic Resource



# System/software Development

Develop data analytic tools

Test important system features

Complete quality control and assurance tasks

# Case Study for Demos & Tutorial

> 30 researcher outreach and training events

> 2000 users

Researcher Workbench  
launched



1 year after Launch

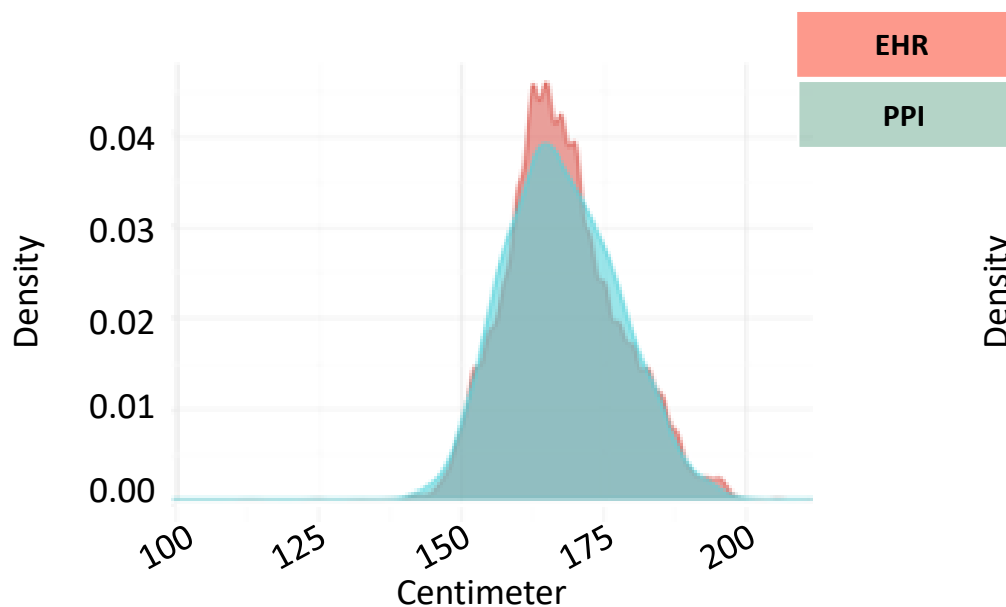
May 2020

May 2021

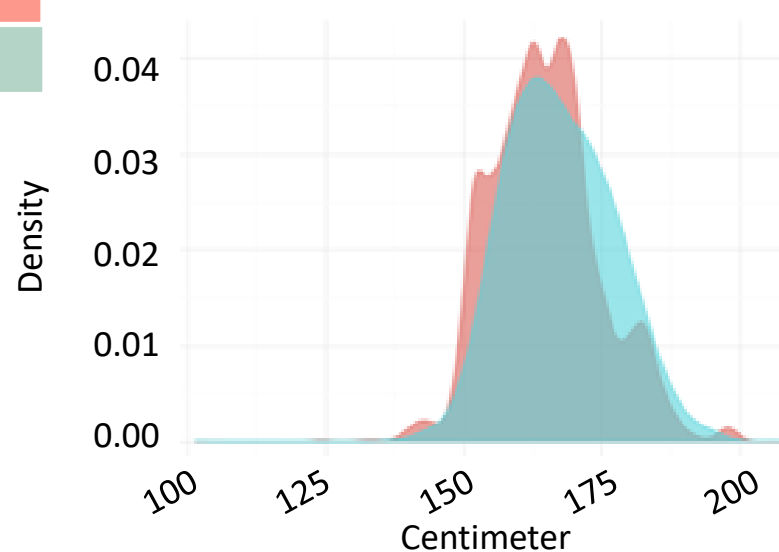
**All of Us**  
RESEARCH PROGRAM



# Real vs Synthetic in the Same Tutorial



Using *real* data in RW

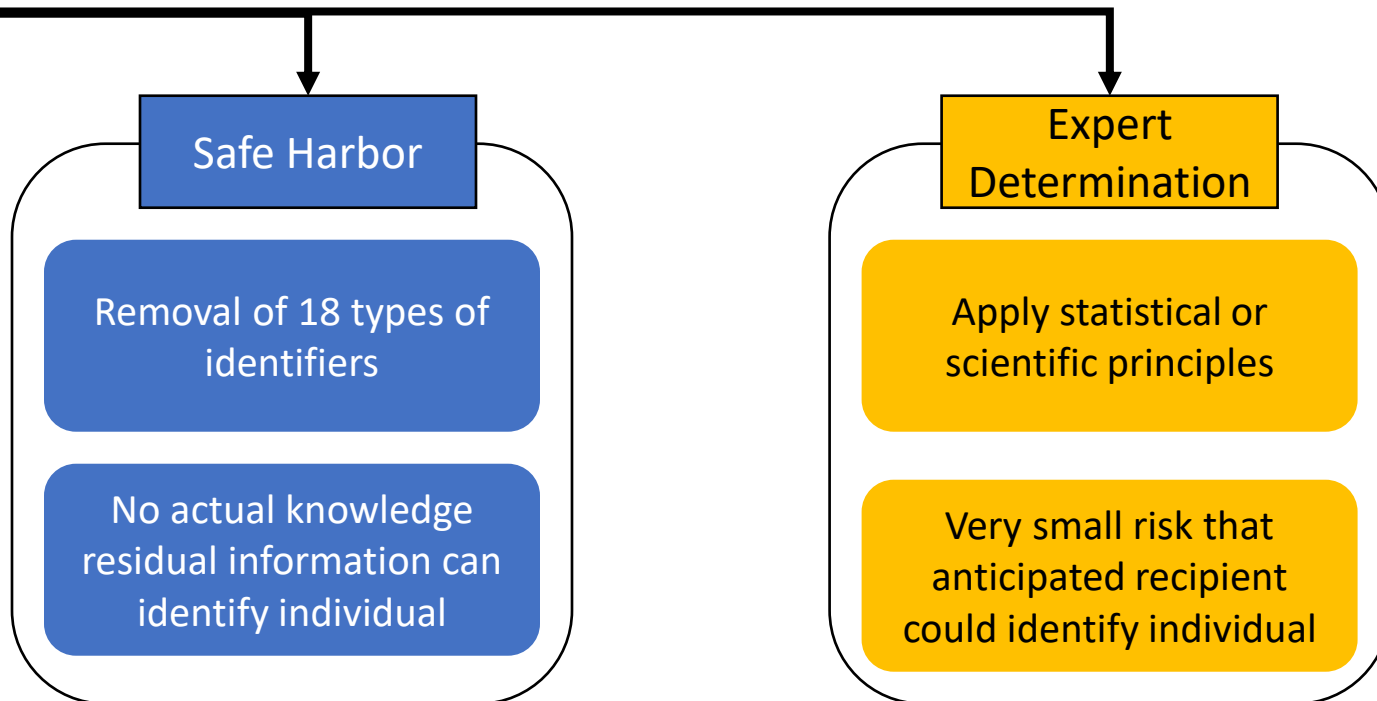


Using *synthetic* data in mirror RW

# Is Synthetic Data “De-identified”?

## According to HIPAA (Privacy Rule):

“information that does not identify an individual and ... no reasonable basis ... information can be used to identify an individual”





What Could Go Wrong?

# AI fake-face generators can be rewound to reveal the real faces they trained on

Researchers are calling into doubt the popular idea that deep-learning models are “black boxes” that reveal nothing about what goes on inside

By Will Douglas Heaven

October 12, 2021

<https://arxiv.org/abs/2107.06304>

## Deep Neural Networks are Surprisingly Reversible: A Baseline for Zero-Shot Inversion

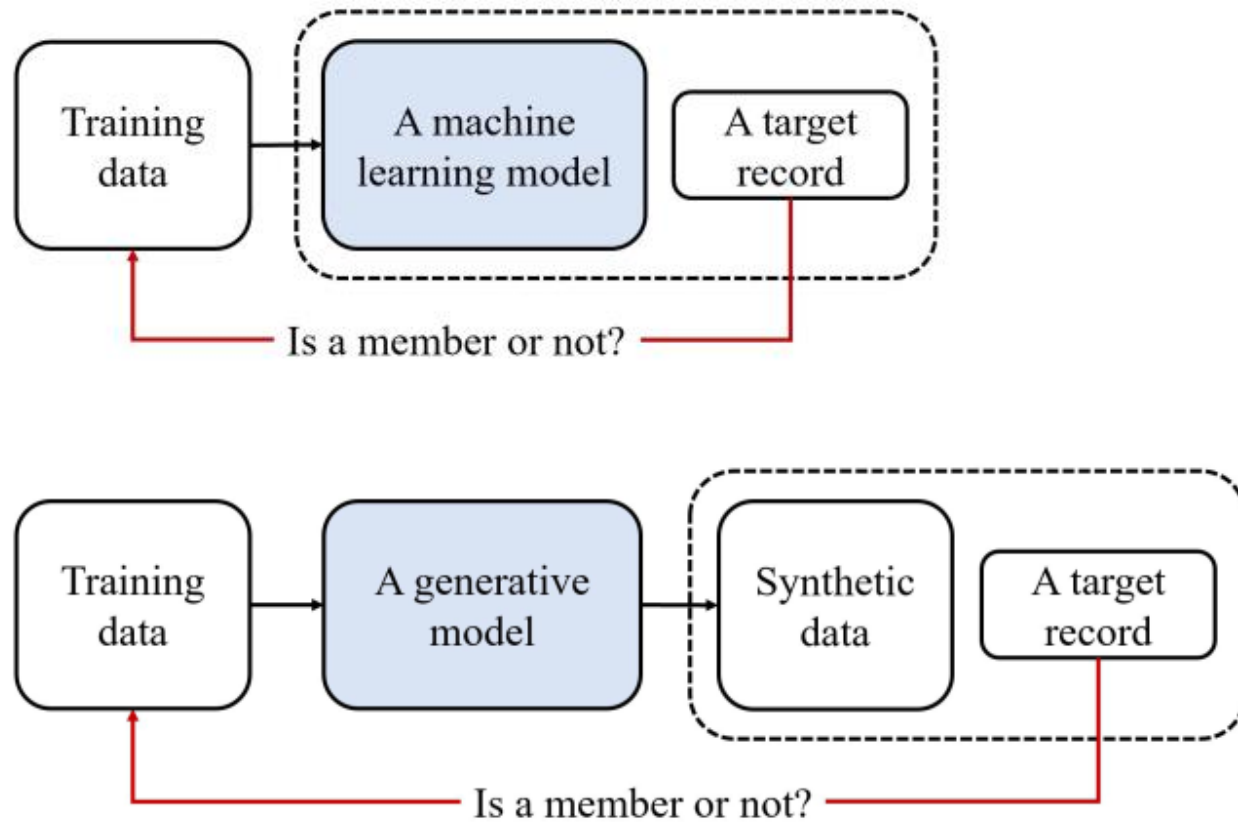
Xin Dong<sup>1,2</sup>; Hongxu Yin<sup>1</sup>; Jose M. Alvarez<sup>1</sup>; Jan Kautz<sup>1</sup>; and Pavlo Molchanov<sup>1</sup>  
<sup>1</sup>NVIDIA, <sup>2</sup>Harvard University  
xindong@g.harvard.edu, {dannyy, josea, pmolchanov, jkautz}@nvidia.com



# A Bunch of Things

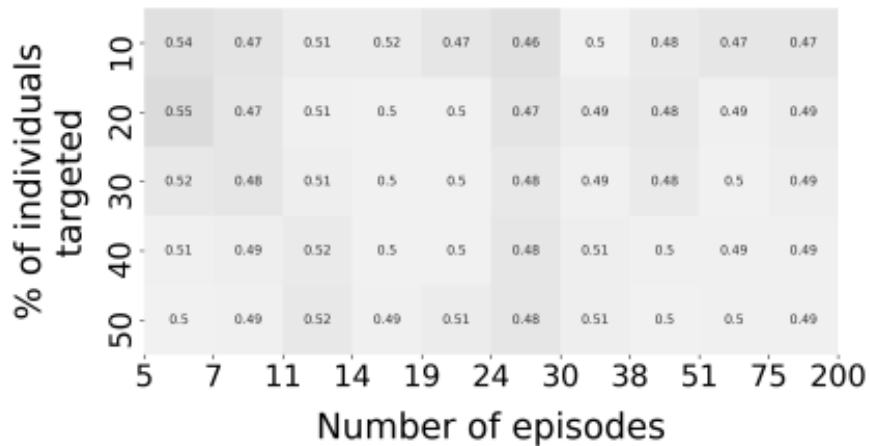
- Mimic
  - Insufficient training data can lead to “mimicking” of original records
- Membership Inference
  - User can test if features of someone they know appear to be in the training data
  - Requires knowing the features in question
- Attribute Inference
  - User can predict features (they don't know) about someone based on features they do know
- Combining Membership and Attribute is where disclosure occurs

# Membership Intrusion

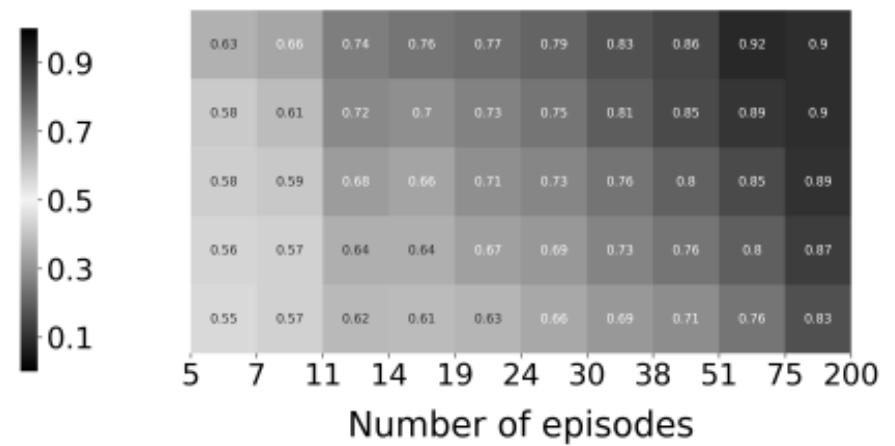


# An Attack on VUMC Data

- 45,000 patients, diagnosis and procedure codes
- Up to 200 visits
- Adversary has 10% “prior” knowledge



Fully Synthetic



Partially Synthetic



# Context Matters ALOT

- Must define the expected capabilities of the recipients of the data
- Privacy assessments should consider the data, as well as how the data was created
- Must consider the recipient's tolerance for errors
- Most consider society's tolerance for intrusion (and claimed intrusion)



# Questions?

b.malin@vanderbilt.edu

Center for Genetic Privacy & Identity in Community Settings

<https://www.vumc.org/getprecise>

Health Data Science Center

<https://www.vumc.org/heads>